Title:        Experimental Design and Comparative Testing of a Hybrid-Cooled
              Computer Cluster Thesis Presentation

Author(s):    Bonnie, Amanda Marie

Intended for: Thesis Defense

Issued:       2015-06-28

# Experimental Design and Comparative Testing of a Hybrid-Cooled Computer Cluster

## Amanda Bonnie

Committee: Payman Zarkesh-Ha, James Plusquellic, and Tarief Elshafiey

June 30, 2015

LA-UR-00-0000

UNCLASSIFIED

# Acknowledgements

- This work was performed using facilities and resources at Los Alamos National Laboratory and was funded by the United States Department of Defense.

# Acknowledgements

- The author thanks her committee members for their review of this work
  - Thank You Dr. Payman, Dr. Plusquellic, and Professor Elshafiey
- The author thanks her direct management for allowing the time to complete this work, along with the support to make it happen.
- The author thanks those at Chilldyne for efforts in water cooling design.
- The author thanks the facilities team and many others at LANL for their added support in facilitating this experimentation
- The author thanks her husband for his continued support through all adventures in life.

# Outline

- Introduction
  - Contributions
- Background
- Theory
- Testbed Overview
  - The Cluster
  - The Water

- Testing Setup
- Results
- Discussion
- Conclusion
- Future Work

# Introduction

1. HPC is growing towards exascale; machine room and/or data center is expanding.

2. Cluster density is growing; more to cool.

3. DOE mandated PUE

4. Total cost of ownership concerns; nearly 30% of a data center electricity bill is spent on cooling

# Contributions

- User-space LDMS dameons

- Deployment, configuration, and support of the TAMIRS cluster at LANL.

- Contract management and installation of hybrid water cooling system.

- Integration of a test suite for monitoring and benchmarking the system for comparative analysis

# Background

- Liquid cooling is NOT novel.

- Cray-2: immersion cooled in Fluorinert in the early 1980s

- NREL: warm water cooling, using waste heat as a main heat source for heating the building

- Sequoia: LLNL reached Top 500 with warm water cooled cluster @ 16 petaflops

- IBM: showed ability to reach 34% increase in processor frequency resulting in 33% increase in performance over same air cooled node

**THEORY**

COSTS, POWER, and JITTER

# Theory – Costs

- Water cooling still costs money BUT…
  - AIR Cooling is expensive too!

- The Data Center of this study has 18 Computer Room Air Conditioning Units (CRACS)
  - 18 consume ~350kW of power
    - Providing 2069kW of cooling capacity
  - At 350kW for the room and $0.1256/kW/h
    - $385,100/year to run the CRAC units (not counting water)

- The full scale water system uses 3 kW of power to cool 200 kW

# Theory – Costs

- AIR:
  - 2069kW Cooling / 350kW Power
  - 5.9 cooling / power
- WATER:
  - 200kW Cooling / 3kW Power
  - 66.66 cooling / power

# Theory – Power

- Fans use power too!

- The nodes for this study have:
  - 6 X Nidec UltraFlow 12VDC, 2.31A fans @ 158CFM

- At full speed that is nearly **166W per node**

- For **20** nodes that **3kW** of power just for fans!

# Theory – Jitter/Noise

- Jitter (noise) can be caused by various means:
  - OS, CPU, and many components

- HPC codes are tightly coupled across the entire job
  - Slowest node, with the slowest core, slows down the entire job
  - Decreases performance of the overall job
  - Increase the job run time

UNCLASSIFIED

# TESTBED OVERVIEW
CLUSTER & INSTALLING THE WATER

# Testbed Overview - Cluster

- TAMIRS:
  - Tiered Active Multi-dimensional Indexed Record Store
- 22 Dell PowerEdge R920 nodes
- 2 Custom SuperMicro nodes
- 4 Dell PowerEdge R720 management nodes
- Purpose: Exploration of next generation tiered storage technologies
- Shared for this study
  - 4 air nodes & 4 water nodes

# Testbed Overview - Cluster

- TAMIRS:
  - Tiered Active Multi-dimensional Indexed Record Store
- 22 Dell PowerEdge R920 nodes
- 2 Custom SuperMicro nodes
- 4 Dell PowerEdge R720 management nodes
- Purpose: Exploration of next generation tiered storage technologies
- Shared for this study
  - 4 air nodes & 4 water nodes

Rack 1　　Rack 2　　Rack 3　　Rack 4

Reserved

# Testbed Overview - Cluster

- Air
  - Node 1 – 4
- Water
  - Node 11 – 14

TAMIRS

I HATE ACRONYMS

CHILLDYNE

NODE 4

NODE 3

NODE 2

NODE 1

NODE 14

NODE 13

NODE 12

NODE 11

CDU 1

CHILLER

CDU 2

# The Install

- Chilldyne Inc,
    - Standalone Chiller Unit
        - Control Inlet Temperature
    - CDU rack
        - CDU X 2 + Vacuum X 2
    - Under floor tubing
    - Above floor tubing
    - Water Block Install

CPU 1

CPU 0

**FINAL DESIGN**

**STOCK**

CPU 1

CPU 0

MEMORY RISER E

MEMORY RISER D

MEMORY RISER B

CPU 1

CPU 0

# TESTING SETUP
TEST CONFIGURATION & MONITORING

# **Testing Setup**

- **<u>Pavillion</u>**: Testing harness in development @ LANL
  - Allows for building a test suite to run the same consistent configurations multiple times
  - Integrated LDMS daemon tool to launch with job
- **<u>HPL</u>**: High Performance LINPACK
- **<u>Systemburn</u>**: Software package from ORNL to create methodical system loads
  - DGEMM: double precision matrix multiplication
  - DSTREAM: double precision floating point vector streaming.
  - PV3: power hungry streaming computational algorithm

# Testing Setup

| Test Name | What is Run | Time To Run [MIN] | X5 |
|:---:|:---:|:---:|:---:|
| HPL | HPL.dat | ~84 | 420 |
| DGEMM | DGEMM_LARGE & DGEMM_SMALL & SLEEP | 105 | 525 |
| DSTREAM | DSTREAM & SLEEP | 45 | 225 |
| PV2 | PV2 & SLEEP | 45 | 225 |
| | **TOTAL TIME:** | **279** | **1395** |

## ~ 24 HOURS OF RAW TESTING

# Testing Setup – Monitoring

- **<u>LDMS</u>**: Lightweight Distributed Metric System
  - Daemon that runs on the nodes sending data to collective source during job run
  - Temperature via lm_sensors
- **<u>PDU</u>**: Power Distribution Units
  - APC AP8641 allow for measurement at each individual plug
- **<u>RAPL</u>**: Running Average Power Limit
  - Intel tool for power capping
  - Allows for power metering at the CPU

UNCLASSIFIED

# RESULTS

## PERFORMANCE, TEMPERATURE, RAPL, & PDU DATA

# Performance HPL

| Cooling Method | Result [GFLOPS] | STDEV | % Improved |
|:---:|:---:|:---:|:---:|
| AIR | 1247 | 9.37 | ---- |
| WATER (65°) | 1257 | 11.04 | 0.80 % |
| WATER (75°) | 1260 | 14.13 | 1.04 % |

# Performance DGEMM

| Cooling Method | MIN. [MFLOPS] | MEAN [MFLOPS] | MAX. [MFLOPS] | % Improved (Mean) |
|---|---|---|---|---|
| AIR | 356.19 | 367.44 | 381.63 | ---- |
| WATER (65°) | <u>371.11</u> | 378.93 | 393.16 | 3.12 % |
| WATER (75°) | <u>366.11</u> | 375.54 | 387.00 | 2.20 % |

~3-4 % higher minimum

# Performance DSTREAM

| Cooling Method | MIN. [MFLOPS] | MEAN [MFLOPS] | MAX. [MFLOPS] | % Improved (Mean) |
|---|---|---|---|---|
| AIR | 354.19 | 360.32 | 366.58 | ---- |
| WATER (65°) | 348.81 | 360.19 | 367.92 | -0.04 % |
| WATER (75°) | 353.94 | 361.81 | 367.47 | 0.45 % |

# Performance PV2

| Cooling Method | MIN. [MTRIPS/s] | MEAN [MTRIPS/s] | MAX. [MTRIPS/s] | % Improved (Mean) |
|---|---|---|---|---|
| AIR | 22.49 | 23.47 | 24.13 | ---- |
| WATER (65°) | 23.92 | 24.09 | 24.18 | 2.66 % |
| WATER (75°) | 23.90 | 24.08 | 24.192 | 2.60 % |

~6% higher minimum

UNCLASSIFIED

# RESULTS
## TEMPERATURE

# Temperature HPL

# Temperature DGEMM

# Temperature DSTREAM

# Temperature PV3
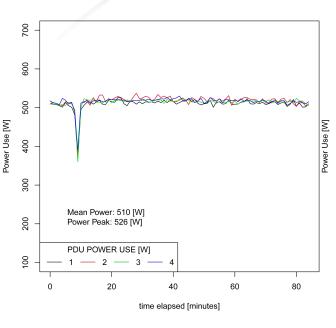
# RESULTS

POWER: RAPL

# RAPL HPL

# RAPL DGEMM

# RAPL DSTREAM

# RAPL PV3

# RESULTS
POWER: PDU

# PDU HPL

# PDU DGEMM

# PDU DSTREAM

# PDU PV3

# DISCUSSION

PERFORMANCE, TEMPERATURE, RAPL & PDU DATA

# Discussion – Performance

- Particular architecture did not benefit from being cooler
  - Power limits could not be turned off to enable longer bursts of turbo clock

- Systemburn showed improved MINIMUM performance values.
  - Increased Min. by **4.19%** and **6.27%** from air to water in DGEMM and PV2 respectively

- At scale this could have a greater effect on overall performance as it seems to have reduced jitter.

# Discussion – Temperature

- Temperature of the core was cooler!
  - 20°C cooler with 65°F water
  - 15°C cooler with 75°F water

- Tighter temperature bands were observed node to node.

- If temperatures were warm enough on even one core, that core could throttle.
  - Keeping it cooler to start, prevents the chance of thermal throttling.

# Discussion – Power: RAPL

- The estimated power use between air and water cooled CPUs was about the same.

- The plots showed a much smoother representation of the package power usage
  - Air tests had lots of jumps
  - Water tests were smooth

- It is not clear what method is used under the covers with RAPL; but from changes in plots, it is expected temperature is considered in some form.

- Perhaps this shows reduced leakage current?

# Discussion – Power: PDU

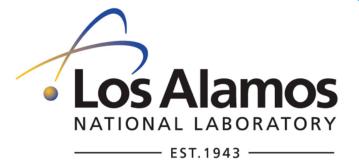- Power data from the PDUs provided proof of reduce power use between air and water cooled nodes.

- Hard to argue:
  - Fans alone or CPUs running cooler? RAPL not so clear.

- BUT! ~30W power savings per node was observed.

- Over the 4 nodes → 120W

- At 2.88kWh/day → $131 saved in a year (for 4 nodes)
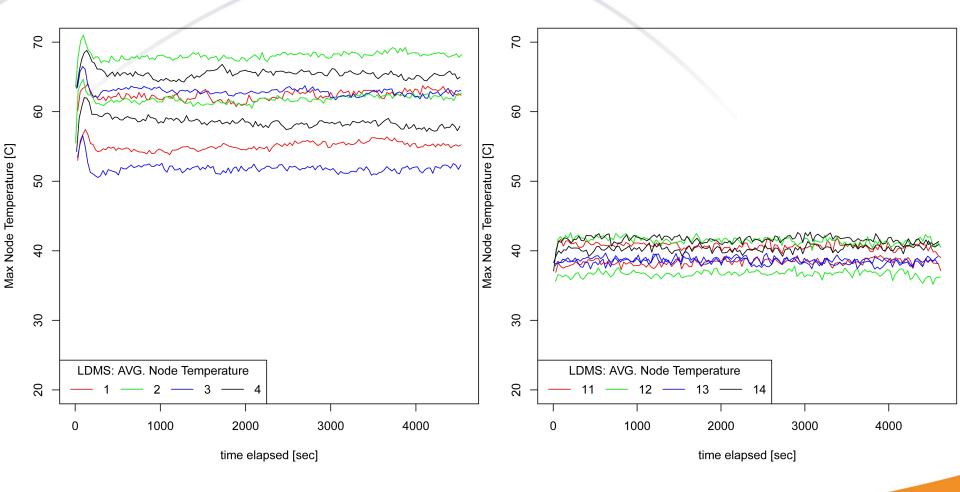
- Scaled up to 20 nodes → 14.4/kWh/day → $658 / year

**BONUS!**

INTERESTING FIND

# BONUS - Data

# **Conclusions**

- Not a major performance gain on average from water alone.

  – Minimum performance improved.

- Idle and load temperatures were significantly reduced.

  – COULD provide greater longevity and better resiliency

- Large cost savings at scale.

  – Factor in reduced CRAC and CHILLER costs!

- RAPL suggests little to no power efficiency gains at chip level.

  – Perhaps reduced leakage current?

UNCLASSIFIED

# Conclusions

- PDU data showed significant power savings.
  - 30W per node → TAMIRS ~600W savings
  - Could significantly scale to clusters of Trinity's size of over 19,000 nodes.
- BONUS: Manufacture layouts can cause some issues.
  - 10°C Temperature difference between CPU0 and CPU1.
  - This could be the difference in a slow core.
  - Widespread temperature band on air cores.
- Lots of potential for power and cooling cost savings.

# Future Work

- **Warm Water Cooling:**
    - Plan to test warmer inlet temperatures.
    - Need more nodes to help maintain the warmer water with dummy loads.
    - Extrapolation from this test suggests 101°F inlet. temperature would be SAFE.
- **Tightly Coupled Applications:**
    - Synthetic benchmarks used were designed to maximize power use and throughput.
    - Plan to test more synchronous dependent workloads.

# Future Work

- **Looking at Scale:**
    - Cluster being a shared resource was not able to be retrofit completely with water cooling.
    - Plan is to get the rest of the 20 compute nodes under water and do more testing.
- **Publish:**
    - SC '15 Power Workshop (August Deadline)
    - Full SC '16 Paper?

# QUESTIONS?

THANK YOU